

SAFETY & RED-TEAMING TAXONOMY (OVERVIEW)

Purpose: Provide a structured, comprehensive taxonomy for evaluating LLM safety, robustness, and policy compliance. Designed for red-teaming, QA, Incident Response Protocols, and governance reviews.

Mapped to NIST AI RMF functions: Govern, Map, Measure, Manage.



- **UNSAFE CONTENT**
(Harmful or Prohibited Outputs)
 - 1.1 **Violence & Physical Harm**
 - Self-harm instructions
 - Violence against others
 - Weapons construction
 - 1.2 **Sexual Content**
 - Content involving minors
 - Non-consensual content
 - 1.3 **Hate & Abuse**
 - Hate speech
 - Harassment & slurs
 - 1.4 **Illegal Activities**
 - Fraud, scams, identity theft
 - Drug manufacturing



- **MODEL & SYSTEM RISKS**
(Technical Vulnerabilities)
 - 2.1 **Prompt Injection**
 - Direct jailbreaks
 - Indirect injection
 - 2.2 **Data Exfiltration**
 - Training data extraction
 - System prompt leakage
 - 2.3 **Hallucinations**
 - Confident falsehoods
 - Fabricated citations
 - 2.4 **Model Bias & Fairness**
 - Stereotyping
 - Discriminatory outputs
 - 2.5 **Privacy & PII Leakage**
 - Re-identification risks

SAFETY & RED-TEAMING TAXONOMY (OVERVIEW)



- **POLICY EVASION & MANIPULATION**

- **3.1 Jailbreak Attempts**

- DAN-style personas
 - Role-play exploits
 - “For research purposes only...”

- **3.2 Obfuscation**

- Encoding (Base64, ROT13)
 - Foreign languages
 - Emojis or spacing tricks

- **3.3 Delegation Attacks**

- Agent-to-agent instruction
 - Multi-step evasion



- **AGENTIC & AUTONOMY RISKS**

- **4.1 Tool Misuse**

- Unauthorised API calls
 - Harmful actions via tools

- **4.2 Excessive Autonomy**

- Running tasks without approval
 - Recursive loops

- **4.3 Safety Bypass via Tools**

- Guardrail circumvention



- **ORGANISATIONAL & OPERATIONAL RISKS**

- **5.1 Compliance Violations**

- UK GDPR, EU AI Act, HIPAA, PCI
 - Lack of transparency

- **5.2 Misalignment**

- Policy contradictions
 - **5.3 Reputational Harm**

- Harmful customer-facing outputs

SAFETY & RED-TEAMING TAXONOMY (APPLICATION)

HOW WE USE THIS TAXONOMY

- This taxonomy serves as the foundational framework for our end-to-end safety lifecycle:
- **Red-Teaming Design:** Defining the scope and specific attack vectors for adversarial testing campaigns.
- **Quality Assurance (QA):** Systematic verification of safety guardrails and policy compliance before deployment.
- **Incident Response Protocols:** Categorising observed failures and near-misses to improve response protocols.
- **Governance Reviews:** Providing a standardised language and metric set for internal and external safety audits.